

"EXPRESS MAIL" Mailing Label No.EL851564935US.....
Date of Deposit. December 13, 2001.....

LOAD BALANCING TECHNIQUE FOR A RESILIENT PACKET RING

CROSS-REFERENCE TO RELATED APPLICATION(S)

5 **[0001]** This application discloses subject matter
related to the subject matter disclosed in commonly
owned, co-pending U.S. Patent Application No.:
_____ (Atty. Docket No. 1285-0078US),
entitled "SIGNALING FOR CONGESTION CONTROL, LOAD
BALANCING, AND FAIRNESS IN A RESILIENT PACKET RING",
10 filed _____ in the name(s) of: Kamakshi
Sridhar.

BACKGROUND OF THE INVENTION

15 Technical Field of the Invention

20 **[0002]** The present invention generally relates to
resilient packet rings ("RPRs") and wavelength division
multiplex RPRs ("WDMRPRs"). More particularly, and not
by way of any limitation, the present invention is
directed to a load balancing technique for use in RPRs
and WDMRPRs.

Description of Related Art

[0003] Resilient Packet Ring ("RPR") is a packet-based technology designed primarily for implementing metropolitan area networks ("MANs"). The basic topology of RPR is a dual counter-rotating ring in which both rings transport both data and control packets. WDMRPR is an RPR with more than one wavelength per ring. RPR technology includes a Medium Access Control ("MAC") mechanism that controls access to the RPR.

[0004] FIG. 1 is an example RPR network 100 consisting of five nodes, respectively designated A-E, and two rings, illustrated in FIG. 1 as an outer ring X and an inner ring Y. Each node A-E is connected to its adjacent nodes via at least two links; one on the inner ring and one on the outer ring. Links on the inner ring are designated a-e, while links on the outer ring are designated α - ϵ . Accordingly, nodes A and B are connected by links α and e, nodes B and C are connected by links β and d, nodes C and D are connected by links γ and c, nodes D and E are connected by links δ and b, and nodes E and A are connected by links ϵ and a.

[0005] Information is carried in one and only one direction on a ring, either clockwise or counter-clockwise. As illustrated in FIG. 1, the outer ring X carries data and control packets in the clockwise direction and the inner ring Y carries data and control packets in the counter-clockwise direction. Features of

RPR include spatial reuse and destination stripping, fast protection, media independent MAC, support for layer 2 devices (e.g., Ethernet switches) and layer 3 devices (e.g., routers). RPR operates as a distributed switch at layer 2. Alternatives to RPR include SONET/SDH rings and Ethernet switches. The former, implemented with routers, is expensive and inefficient for dynamic traffic patterns; however, it provides efficient protection. The latter offers no protection bandwidth or fairness, but does provide efficient bandwidth utilization. Moreover, in contrast to layer 3 devices, RPR has ring awareness. RPR, therefore, is intended to embody the best features of both SONET and Ethernet switches, that is, bandwidth efficiency, fairness, and fast protection. To realize these features, the RPR should provide a basic ring access mechanism, additional mechanisms to provide satisfactory performance for traffic transiting the ring, and appropriate protection features.

[0006] Several RPR terms and concepts will now be defined. Referring again to FIG. 1, assume that traffic destined for node D enters the ring at node A and passes through nodes B and C. Node A is referred to as the "ingress node", node D is referred to as the "egress node", and nodes B and C are referred to as "intermediate nodes." As previously noted, the fiber connecting two nodes is referred to as a "link". In WDMRPR, a wavelength between two nodes is referred to as a link. From the perspective of nodes B and C, traffic from node

A is "ring traffic". Traffic entering node B or C is referred to as "tributary traffic," with respect to node A traffic.

[0007] At the ingress node, as soon as a flow is admitted to the RPR, an RPR header is created and attached to every packet within the flow. The RPR header contains a Class of Service ("CoS") marking, a preemption bit, a ring identifier, and a destination MAC address. Each intermediate node will examine the incoming packet and perform an address lookup. If the packet is destined for the current node, it is taken off the ring at that node. If it is destined for some other node, it is placed back on to the ring.

[0008] One problem inherent in RPR architecture is the problem of load imbalance between two or more rings or wavelengths. Hereinafter, unless otherwise explicitly specified, the term "rings" will be deemed to encompass both rings and wavelengths. Load imbalance results from excessive traffic on one ring, as compared to little traffic on another ring. In an RPR with two rings, there are only two options for sending traffic from one node to another node, either on one ring or on the other. In such cases, traffic is typically sent on the ring with the fewest hops. This may result in a situation where a lot of traffic travels from one node to another on only one ring, because the ring has fewer hops as compared to the other ring.

[0009] Additionally, because RPR employs destination stripping, certain segments on one ring are typically more heavily loaded as compared to other segments on the same ring. Also, certain segments on one ring may be much more heavily loaded as compared to segments on the other ring.

[0010] As previously indicated, load imbalance may be defined as the concurrent overutilization of capacity on one or more links in the RPR and the underutilization of capacity of other links in the RPR. Load imbalance can be illustrated as follows. Referring again to FIG. 1, traffic going from node A to node B can travel either on the outer ring, via link α , or on the inner ring via links a, b, c, and d. For obvious reasons, traffic will tend to take the link α , since it is the more direct route. In the absence of any other traffic on the RPR, the link α will be overutilized, while the other links will be underutilized. The solution is to "load balance" the traffic between the links so that packets of a given class will experience approximately the same delay regardless of what the shortest path is.

[0011] Thus, traffic going from one node to another experiences varying amounts of delay depending on how many and which heavily loaded segments it traverses. Since traffic performance is judged by total end-to-end delay, the traffic performance (delay, dropped packets, etc.) on one ring may be significantly worse than the

traffic performance on the other ring for a single class of service. This results in congestion on one ring versus unused capacity on the other ring and therefore requires load balancing.

5 **[0012]** Load balancing in RPR involves two functions: (a) diverting incoming flows from more heavily loaded rings to alternate paths on other rings, where they will likely experience more hops, but better performance, due to less traffic; and (b) improving the performance of
10 existing traffic on the heavily loaded ring. Perfect load balancing in RPR, for all classes, between all pairs of nodes, is challenging due to the ring structure and destination stripping features of RPR. Load balancing performed for one pair of nodes for one class may result
15 in traffic redistribution that might necessitate load balancing between another pair of nodes for the same or different traffic class. Dynamic traffic patterns and time-dependent flow arrivals and departures also complicate load balancing in RPR.

20 **[0013]** Prior art for load balancing traffic in multiple rings (or wavelengths) of RPR is currently not known to exist. Prior art for changing quality of service ("QoS") parameters of various classes relative to one another is inadequate to help divert traffic to
25 alternate paths on multiple rings. Prior art to change the token bucket ("TB") parameters, including peak data rate and maximum burst size, for aggregate traffic flows does not allow the behavior of traffic classes to be

changed relative to one another. While it may be adequate to reduce congestion in a network, it is inadequate to implement load balancing between rings in an RPR. Finally, all known prior art is restricted to only TB parameters and does not consider the interaction between the bandwidth broker ("BB") parameters, which include allocated bandwidth, available bandwidth, and used bandwidth, and the QoS parameters, which include the TB parameters as well as the class based queuing ("CBQ") parameters.

SUMMARY OF THE INVENTION

[0014] Accordingly, the present invention advantageously provides a technique for load balancing in an RPR. In one embodiment, the present invention comprises implementing on every node a QoS/BB monitor to achieve load balancing in an RPR. The principles of the present invention are applicable to dual-ring and multi-ring RPR, as well as wavelength division multiplexing ("WDM") and dense wavelength division multiplexing ("DWDM") RPR. The term "RPR" will be deemed to include all of the foregoing types of RPR, as well as other types of RPR to which the principles of this invention may apply. The QoS/BB monitor on a node is common to all rings in the RPR. It has knowledge of traffic performance for each class on each ring of the RPR, which knowledge is obtained through periodic measurements or in response to failure events. This allows the monitor to

vary the QoS parameters on a node, for one or more rings,
for a particular traffic class, to achieve load
balancing. Likewise, the QoS/BB monitor can signal to
the BB entity at higher layers to vary the BB parameters
on a node, for one or more rings, for a particular class,
to achieve load balancing.

[0015] Changing BB parameters affects call admission
of new traffic flows. Changing QoS parameters affects
performance of traffic already on the ring. BB parameter
variations result in coarse variation of traffic
performance, while QoS parameter variations are likely to
result in finer variation of traffic performance.
Through the combination of varying QoS parameters and BB
parameters, the QoS/BB monitor can throttle traffic on
one ring and permit more traffic on another ring, for a
particular traffic class, to achieve load balancing. By
signaling to the BB entity to change its parameters, the
ability to divert traffic at the node ingress from one
ring to another ring may be achieved.

[0016] The key benefit of having one entity that
controls both the QoS parameters and the BB parameters is
that the traffic behavior can be modified appropriately,
in terms of direction and servicing rate, to ultimately
achieve load balancing.

[0017] In one aspect, the invention comprises a method
of implementing load balancing in an RPR network
comprising a plurality of nodes and first and second

5

20

bandwidth utilization on the less heavily loaded one of the rings.

BRIEF DESCRIPTION OF THE DRAWINGS

5 **[0020]** A more complete understanding of the present invention may be had by reference to the following Detailed Description when taken in conjunction with the accompanying drawings wherein:

10 **[0021]** FIG. 1 depicts an exemplary RPR network arrangement in which teachings of the present invention may be advantageously practiced;

[0022] FIG. 2 (Prior Art) depicts a block diagram of an exemplary embodiment of a single node of an RPR network according to teachings of the prior art;

15 **[0023]** FIG. 3 depicts a block diagram of an exemplary and presently preferred embodiment of a single node of an RPR network in which teachings of the present invention may be advantageously practiced; and

20 **[0024]** FIG. 4 depicts a flowchart of a load balancing technique for implementation in an RPR network in accordance with the teachings of the present invention.

DETAILED DESCRIPTION OF THE DRAWINGS

25 **[0025]** In the drawings, like or similar elements are designated with identical reference numerals throughout the several views thereof, and the various elements depicted are not necessarily drawn to scale. Referring

again to FIG. 1, depicted therein is an exemplary RPR network 100 in which the teachings of the present invention may be advantageously practiced. It should be recognized that, although the RPR network 100 comprises five nodes A-E, in general, there can be more or fewer nodes on the network 100. Moreover, although each of the nodes A-E are illustrated as being connected to its adjacent nodes by only two links, it should be recognized that the principles of the invention described herein can be applied to an RPR comprising more than two rings, and therefore more than two links between adjacent nodes. The teachings of the present invention are also applicable to WDM and DWDM RPR networks.

[0026] Consider a large volume of traffic going from node A to node C. Two choices exist: (a) the path including links α and β (path α - β) via the outer (clockwise) ring; and (b) the path including links a, b, and c (path a-b-c) via the inner (counter-clockwise) ring. Typical path selection protocols will result in the selection of the first path (path α - β), as it has the fewest number of hops. New traffic entering node A destined for node C must be diverted to path a-b-c for better performance. This is achieved through the load balancing method of the present invention as described herein. Additionally, performance of existing traffic on the path α - β will be improved as part of the load balancing method of the present invention.

5 [0027] Consider now new traffic entering node B destined for node C. The natural path is via link β , which provides a direct path between nodes B and C; however, this link is heavily loaded with traffic from node A destined for node C. Therefore, it may be better to send new incoming traffic via the path comprising links e, a, b, c (path e-a-b-c).

10 [0028] To account for traffic variations, load balancing is performed at periodic time intervals on a per class basis. For discussion, it will be assumed that traffic flows are classified into three classes of service, including Class 1 (highest priority), Class 2 (medium priority), and Class 3 (lowest priority); however, it will be recognized that the teaching of the present invention may be applied to any number of classes of service. Additionally, although for the sake of simplicity, the load balancing technique of the present invention will be illustrated as applied to Class 1, it will be recognized that it may be applied to any class of service.

15 [0029] FIG. 2 is a block diagram of an existing QoS implementation within a single node 200. It will be assumed for the sake of example that the node 200 represents node A of the RPR 100 (FIG. 1). In general, as illustrated in FIG. 2, packets entering the node 200 are classified by a packet classifier 222 into one of the three classes. Packets of each class are shaped by a

respective token bucket policer 224 depending on various token bucket ("TB") parameters (e.g., peak data rate and maximum burst size). Packets may be dropped if the traffic rate is excessive. Once shaped, packets are put
5 into respective class based queues ("CBQs") 226 and drained from the CBQs 226 depending on the CBQ parameter specification of how many packets of each class can be serviced by a scheduler 228 at each service rotation. The CBQ and TB parameters together comprise the quality
10 of service ("QoS") parameters.

[0030] A bandwidth broker ("BB") 230 located at a higher layer (e.g., layer 3) is responsible for bandwidth allocation to each class for a particular link. The BB parameters per class include allocated bandwidth, used
15 bandwidth, and available bandwidth. The allocated bandwidth for a class refers to the total bandwidth available for that class on a link. Of this, only a portion of the bandwidth is used by the various flows. The difference between allocated bandwidth and used
20 bandwidth is the available bandwidth for that class on that link. The sum of the allocated bandwidths for all classes equals the useful link capacity.

[0031] In a top portion 200a of the node 200, the QoS parameters are applied to traffic going to link α between
25 nodes A and B and in a bottom portion 200b of the node 200, the QoS parameters are applied to traffic going to link a between nodes A and E. Likewise for the BB

parameters, even though the BB 230 is shown as a single entity for simplicity.

[0032] FIG. 3 is a block diagram of a node 300 embodying features of the present invention for implementing a load balancing technique. The node 300 is identical in all respects to the node 200 (FIG. 2), except that it includes a QoS/BB monitor 302. For purposes of clarity, only interaction of the monitor 302 with Class 1 is illustrated in FIG. 3; however, it should be recognized that, in practice, the monitor 302 is common to all rings in the RPR and monitors the traffic performance for each class on each ring at the node 300. With this knowledge, it can change the QoS parameters of each class on each link one at a time, for a particular traffic class. Likewise for the BB parameters.

[0033] Traffic performance (e.g., delay, jitter, dropped packets) is sensed by the QoS/BB monitor 302 at periodic time intervals by an in- or out-of-band signaling mechanism, such as that described in commonly-assigned U.S. Patent Application Serial No. _____ (Atty. Docket No. 1285-0078US), filed _____, entitled SIGNALING FOR CONGESTION CONTROL, LOAD BALANCING, AND FAIRNESS IN A RESILIENT PACKET RING, which is hereby incorporated by reference in its entirety. Changes to the QoS and BB parameters are made after load imbalance is detected, as will be described in greater detail below.

[0034] The load balancing technique will now be described with reference to FIGs. 1,3 and 4. Load balancing is initiated either after a failure event or at periodic time intervals T_k . The following variables, as defined below, will be used to illustrate the load balancing technique of the present invention. In connection with a variable, subscript 1 denotes Class 1, subscript 2 denotes Class 2, and subscript 3 denotes Class 3. The load balancing technique is described here with reference to node A. The same technique is applicable to every other node. Notation is defined below only for Class 1.

15	$B_{1\alpha}, B_{1a}$	# bytes in CBQ drained off each rotation for traffic going to link α and link a, respectively
	$T_{pl\alpha}, T_{pla}$	TB peak rate allowable for traffic going to link α and link a, respectively
20	$T_{bl\alpha}, T_{bla}$	maximum TB burst size allowable for traffic going to link α and link a, respectively
	T_k	time interval at which traffic measurements are made to detect load imbalance
	$T_{all\alpha}, T_{alla}$	allocated bandwidth for links α and a, respectively
25	$T_{avl\alpha}, T_{avla}$	available bandwidth for links α and a, respectively
	$T_{ul\alpha}, T_{ula}$	used bandwidth for links α and a, respectively
30	$d_{1\alpha}, d_{1a}$	delays of the marked packets on links α and a, respectively, for Class 1 obtained at specific time intervals T_k or when there is a failure event
35	$dr_{1\alpha}, dr_{1a}$	# of dropped packets on link α and link a, respectively, for Class 1 during the time interval T_k

$j_{1\alpha}$, j_{1a} jitter of packets on link α and link a , respectively, for Class 1 (to determine these variables, two successive packets for the class need to be marked)

5

[0035] CBQ parameters for link α are $B_{1\alpha}$, $T_{p1\alpha}$, and $T_{b1\alpha}$, and CBQ parameters for link a are B_{1a} , T_{p1a} , and T_{b1a} . BB parameters for link α are $T_{all\alpha}$, $T_{av1\alpha}$, and $T_{u1\alpha}$. BB parameters for link a are T_{alla} , T_{av1a} , and T_{ula} . It should be noted that available bandwidth and used bandwidth are obtained from the BB entity 230 at higher layers made available through an RPR control channel. Example of constraints to which traffic in class 1 must conform are defined as follows:

15

$$\begin{aligned} dr_{1(\alpha \text{ or } a)} &\leq DR_1 \\ d_{1(\alpha \text{ or } a)} + j_{1(\alpha \text{ or } a)} &\leq D_{max1} \end{aligned}$$

20

where DR_1 is the maximum acceptable number of dropped packets for Class 1 and D_{max1} is the maximum acceptable delay for Class 1. Similar constraints are defined for other classes.

25

[0036] FIG. 4 is a flowchart of the operation of a load balancing method for RPR in accordance with teachings of the present invention. Execution begins in step 400 at a time T . In step 400, load imbalance detection is performed and a determination is made whether a load imbalance has been detected. Load imbalance detection at each node is skewed with respect to other nodes at time intervals of m , where $m < T_k$. Detecting load imbalance may be accomplished in any number of ways, including, for example, measurements of

30

delays and/or dropped packets. One method of detecting load imbalance is to measure and compare delays for test packets sent from node A to a node midway, for example, node C, on each of the two rings. Otherwise, some other
5 node destination may be picked to test for load imbalance. If the delay and/or number of dropped packets is significantly higher on one ring, e.g., ring X, than the other ring, e.g., ring Y, for a given pair of nodes, then it could trigger the load balancing operation, as
10 described below. Alternatively the constraints for each of the rings X and Y could be evaluated to determine a trigger for the load balancing operation. The threshold for triggering the load balancing operation should be sufficiently low so that the QoS/BB monitor can be
15 sensitive to dynamic traffic patterns, but sufficiently high to prevent too frequent load balancing from occurring.

[0037] If load imbalance is not detected in step 400, step 400 is repeated again at the next time interval T_k
20 (step 401). If load imbalance is detected in step 400 between rings X and Y, then load imbalance correction must be performed at node A as described in steps 402-412. Since link α is overutilized and link a is underutilized, the BB parameters must be changed so that
25 traffic is diverted at node A from link α to link a (i.e., from ring X to ring Y). In step 402, responsive to a decision to implement load balancing, the QoS/BB monitor 302 at node A signals the BB entity 230 to change

the BB parameters so that more bandwidth is available for new incoming Class 1 flows on ring Y and less bandwidth is available for new incoming traffic for Class 1 on ring X. Note that except for the interaction with the QoS/BB monitor 302, the two rings operate independently of one another. The QoS/BB monitor 302 signals the BB monitor 230 to decrease $T_{a11\alpha}$ and increase T_{a11a} . This will effectively cause a Call Admission Control (CAC) entity (not shown), either co-located with or separate from the BB entity 230 at node A, to admit more traffic flows on ring Y rather than on ring X, due to the increase in available bandwidth for Class 1 traffic on link a and the concurrent decrease in available bandwidth for Class 1 traffic on link α .

[0038] Next, the QoS parameters for Class 1 on ring X and for Class 1 on ring Y must be changed so as to improve traffic performance on ring X while increasing bandwidth utilization on ring Y. From the perspective of node A, the QoS parameters for Class 1 must be changed so as to improve traffic performance on link α while increasing bandwidth utilization on link a. To let more traffic on link a at node A, the TB peak rate T_{p1a} must be increased. To restrict the traffic on link α at node A, the TB peak rate $T_{p1\alpha}$ must be decreased. The foregoing is accomplished as described below with respect to steps 404-414. In step 404, $T_{p1\alpha}$ is set to $T_{u1\alpha}$. There is an interaction between the BB parameters and the QoS parameters. By signaling changes to the BB entity 230 to

decrease $T_{all\alpha}$ (step 402), fewer flows will be admitted on link α on ring X; thus, $T_{ul\alpha}$ will decrease over a period of time, or not increase very much. Since $T_{p1\alpha}$ is set to $T_{ul\alpha}$ a decrease in $T_{ul\alpha}$ will result in a decrease in $T_{p1\alpha}$. Thus existing traffic for Class 1 will enter link α at node A at a lower peak rate and hence improve the performance of existing traffic on link α for Class 1. In step 406, the TB parameters for Classes 2 and 3 on link α on ring X ($T_{p2\alpha}$, $T_{b2\alpha}$, $T_{p3\alpha}$, and $T_{b3\alpha}$) are reduced. The reduction of these parameters relative to Class 1 traffic helps improve the performance of existing Class 1 traffic on link α relative to other classes on link α on ring X. If load imbalance is not subsequently detected, the Class 2 and Class 3 TB parameters can be increased again.

[0039] In step 408, T_{p1a} is set to T_{u1a} . Again, there is an interaction between the BB parameters and the QoS parameters. By signaling changes to the BB entity 230 to increase T_{all1a} (step 402), more flows will be admitted on link a on ring Y, and thus T_{u1a} will increase. Since T_{p1a} is set to T_{u1a} , an increase in T_{u1a} will result in an increase in T_{p1a} . This will allow more traffic to enter link a at node A and hence increase the bandwidth utilization of Class 1 traffic on link a.

[0040] In step 412, $B_{1\alpha}$ is increased, while $B_{2\alpha}$ and $B_{3\alpha}$ are kept fixed. This preferentially permits more Class 1 traffic to be serviced on link α on ring X at each service rotation of the scheduler 228, relative to that

for Classes 2 and 3, and thus helps improve the performance of Class 1 traffic. In step 418, B_{1a} is increased, while B_{2a} and B_{3a} are kept fixed. This preferentially permits more Class 1 traffic to be serviced on link a at each service rotation of the scheduler 228 relative to that for classes 2 and 3 and thus helps increase bandwidth utilization of Class 1 traffic on link a. The actual amount of increase may be performed in powers of two, for example, or using some other algorithm. Again, if load imbalance is not subsequently detected, the various Class 1 CBQ parameters can be changed on each of the two rings X and Y.

[0041] Changing the QoS parameters affects the performance of traffic already on the ring. Changing the BB parameters affects call admission for new traffic flows. Note that since these parameters can be changed at periodic time intervals or when there is a failure event, there is an opportunity to correct an incorrect selection of parameters at the next time interval T_k using new measurement data. Therefore, the actual amount by which the QoS/BB parameters are increased or decreased affects only the time it takes to accomplish load balancing, but does not alter the objectives and results of load balancing.

[0042] Other options may also be implemented on the QoS/BB monitor. Various load imbalance detection algorithms can be selected. The selection of the measurement time interval is also left to the operator.

The shorter the timer interval T_k , the more measurements that are available and thus, a better judgment of load imbalance can be made at the expense of computation time. This described mechanism for load balancing can easily be extended to RPRs comprising more than two rings or wavelengths.

[0043] Based upon the foregoing Detailed Description, it should be readily apparent that the present invention advantageously provides an innovative and efficient solution for correcting load imbalance in an RPR. In particular, because the QoS/BB monitor is common to all rings in an RPR, load balancing can be implemented easily, with minimal computation. The load balancing algorithm (FIG. 4) works by exercising control of traffic ingress via the BB parameters, and of traffic already on the RPR via the QoS parameters. The load balancing technique described herein allows correction of load imbalance on a class-by-class basis. Incorrect guesses for the QoS parameters and/or BB parameters can be corrected at the next time interval when new measurements are available. Finally, coordination between the QoS/BB monitors of various nodes can be implemented for other load balancing, congestion control, and fairness algorithms.

[0044] It is believed that the operation and construction of the present invention will be apparent from the foregoing Detailed Description. While the exemplary embodiments of the invention shown and

described have been characterized as being preferred, it
should be readily understood that various changes and
modifications could be made therein without departing
from the scope of the present invention as set forth in
5 the following claims.